

AD-A104 860 RICE UNIV HOUSTON TEX DEPT OF MATHEMATICAL SCIENCES F/G 12/1
AN ALGORITHM FOR NONPARAMETRIC DENSITY ESTIMATION, (U)
MAY 76 D W SCOTT, R A TAPIA, J R THOMPSON E-(40-1)-5046

RICE UNIV HOUSTON TEX DEPT OF MATHEMATICAL SCIENCES
AN ALGORITHM FOR NONPARAMETRIC DENSITY ESTIMATION; (U)
MAY 76 D W SCOTT; R A TAPIA; J R THOMPSON E-140

F/G 12/1

UNCLASSIFIED

E-(40-1)-5046

NL

116

END

DATE _____

THE MFC

10 x

DTIC

LEVEL

For: Computer Science and Statistics: Ninth Annual Symposium on the Interface.

May 1976

DTIC
ELECTE

OCT 1 1981

2

AD A104800

AN ALGORITHM FOR NONPARAMETRIC DENSITY ESTIMATION

David W. Scott
Dept. of Biostatistics
Baylor College of Medicine
Houston, Texas 77030

Richard A. Tapia
Dept. of Mathematical
Sciences
Rice University
Houston, Texas 77001

James R. Thompson
Dept. of Mathematical
Sciences
Rice University
Houston, Texas 77001

ABSTRACT

A numerical algorithm is given for implementing a nonparametric maximum penalized likelihood estimator similar to those proposed by Good and Gaskins and those proposed by de Montricher, Tapia and Thompson. It is shown how the resulting nonlinear constrained optimization problem may be effectively solved by using Tapia's approach to Newton's method for constrained problems.

1. Introduction. de Montricher, Tapia and Thompson demonstrated that the standard histogram was an unstable maximum likelihood density estimator and considered maximum penalized likelihood estimators similar to those previously considered by Good and Gaskins (1971). Specifically suppose we are given the random sample $x_1, \dots, x_n \in (a, b)$. Let $H_0^1(a, b)$ consist of the functions f defined on (a, b) with the property that $f(a) = f(b) = 0$ and f' is a member of $L^2(a, b)$. Estimate the density function which gave rise to the random sample x_1, \dots, x_n by the solution of the constrained optimization problem

$$(1.1) \max L(f); f \in H_0^1(a, b), f \geq 0 \text{ and}$$

$$\int_a^b f(x) dx = 1,$$

where

$$(1.2) L(f) = \prod_{i=1}^n f(x_i) \exp(-\alpha \int_a^b |f'(x)|^2 dx), (\alpha > 0).$$

The functional L in (1.2) is called the penalized likelihood and the solution of (1.1) is called the maximum penalized likelihood estimator based on the random sample x_1, \dots, x_n . de Montricher, Tapia and Thompson (1975) proved that problem (1.1) has a unique solution and is a monospline of degree two,

i.e., a polynomial of degree two plus a spline of degree one. We now give a numerical algorithm for approximating this monospline.

2. The Discrete Problem. For given n , consider the mesh t_0, \dots, t_{n+1} where $t_i = a + ih$, $i = 0, \dots, n+1$ with $h = (b-a)/(n+1)$. Let H_0^1 denote the vector space of all continuous piecewise linear functions which have knots at t_1, \dots, t_n and vanish at a and b . For $p \in H_0^1$ let $y_i = p(t_i)$, $i = 0, \dots, n+1$. Then $y_0 = y_{n+1} = 0$ and

$$(2.1) p(x) \geq 0 \Leftrightarrow y_i \geq 0, \quad i = 1, \dots, n$$

$$(2.2) \int_a^b p(x) dx = h \sum_{i=0}^n y_i$$

$$(2.3) \int_a^b p'(x)^2 dx = \frac{1}{h} \sum_{i=0}^n (y_{i+1} - y_i)^2.$$

Let

$$(2.4) v_1 = \# \text{ of } x_i \text{ in } [a, t_1 + \frac{h}{2})$$

$$(2.5) v_i = \# \text{ of } x_i \text{ in } [t_{i-1} + \frac{h}{2}, t_i + \frac{h}{2}), \quad i = 2, \dots, n-1$$

$$(2.6) v_n = \# \text{ of } x_i \text{ in } [t_{n-1} + \frac{h}{2}, b).$$

We shall assume that we have enough data so that $v_i > 0 \forall i$. Our finite dimensional

This work was supported in part by ONR grant ONR-042-283 and ERDA contract E-(48-1)-5646.

81 9 30 078

This document has been approved for public release and sale; its distribution is unlimited.

DTIC FILE COPY

4-573

approximation to problem (1.1) is

$$(2.7) \max \hat{L}(y); y_i \geq 0 \forall i \text{ and } \sum_{i=1}^n y_i = h^{-1}$$

$$(2.8) \hat{L}(y) = \prod_{i=1}^n y_i \exp(-\alpha h^{-1} \sum_{i=1}^n (y_{i+1} - y_i)^2).$$

Clearly (2.7) is a constrained optimization problem in R^n .

Proposition 2.1. The constraints $y_i \geq 0$ of problem (2.7) are not active at the solution.

Proof. If $y^* = (hN)^{-1}(1, \dots, 1)$, then y^* satisfies all the constraints of problem (2.7) and $\hat{L}(y^*) > 0$. Moreover, if $y = (y_1, \dots, y_n)$ is such that $y_i = 0$ for some i , then $\hat{L}(y) = 0$. This proves the proposition.

It follows that we can obtain the solution of problem (2.7) by solving

$$(2.9) \min(-\log \hat{L}(y)); \sum_{i=1}^n y_i = h^{-1}$$

where from (2.8) we see that

$$(2.10) -\log(\hat{L}(y)) = -\sum_{i=1}^n v_i \log(y_i) + \alpha h^{-1} \sum_{i=1}^n (y_{i+1} - y_i)^2.$$

3. The Algorithm. The Lagrangian for problem (2.9) is

$$(3.1) \mathcal{L}(y, \lambda) = -\sum_{i=1}^n v_i \log(y_i) + \alpha h^{-1} \sum_{i=0}^n (y_{i+1} - y_i)^2 + \lambda (\sum_{i=1}^n y_i - h^{-1}).$$

The gradient of the Lagrangian is

$$(3.2) \nabla \mathcal{L}(y, \lambda) = (\dots, -v_i y_i^{-1} + 2\alpha h^{-1}(-y_{i+1} + 2y_i - y_{i-1}) + \lambda, \dots)^T$$

and the Hessian of the Lagrangian is the diagonally dominant tridiagonal matrix

$$(3.3) \nabla^2 \mathcal{L}(y, \lambda) = \begin{pmatrix} d_0^1 & d_1 & & & \\ d_{-1} & d_0^2 & d_1 & & \\ & d_{-1} & d_0^{n-1} & d_1 & \\ & & d_{-1} & d_0^n & \\ & & & d_{-1} & d_0^n \end{pmatrix}$$

where $d_{-1} = d_1 = -2\alpha h^{-1}$ and

$$d_0^1 = 4\alpha h^{-1} + v_1 (y_1^{-1})^2.$$

It therefore follows that Tapia's (1974), (1976) approach to Newton's method for constrained problems is a natural one for this problem and the operation count will be

$O(n)$ per iteration instead of the usual $O(n^3)$ expected from Newton's method.

Let $g(y) = \sum_{i=1}^n y_i - h^{-1}$. Then

$U = \nabla g(y) = (1, \dots, 1)^T$. We use \langle, \rangle to denote the inner product in R^n .

The Newton-like Algorithm.

Step 1. Determine $\alpha > 0$, $\epsilon > 0$, y^0, λ^0 and set $k := 0$.

Step 2. Calculate $\lambda^{k+1} = \langle U, \nabla^2 \mathcal{L}(y^k, \lambda^k)^{-1} U \rangle^{-1} \langle U, \nabla \mathcal{L}(y^k, \lambda^k)^{-1} U \rangle$ and $y^{k+1} = y^k - \nabla^2 \mathcal{L}(y^k, \lambda^k)^{-1} \nabla \mathcal{L}(y^k, \lambda^{k+1})$

Step 3. If $\|\nabla \mathcal{L}(y^k, \lambda^k)\| \leq \epsilon$, then stop. If not, then set $k := k+1$ and go to Step 2.

Initialization values could be

$$\epsilon = 10^{-4} \\ \alpha = 5.0 \\ y^0 = (nh)^{-1}(1, \dots, 1) \\ \text{and} \\ \lambda^0 = -2(nh)^{-1}(y_1 + y_n) + \sum_{i=1}^n v_i (ny_i)^{-1}.$$

This λ^0 is given by the projection formula in Tapia (1976).

For a complete description of this algorithm and related quasi-Newton methods for constrained optimization the reader is referred to Tapia (1976).

Proposition 3.1. The above algorithm is locally quadratically convergent and requires only $O(n)$ operations per iteration.

4. Some Numerical Examples. Although for reasons of conciseness it was appropriate to develop the above discrete maximum penalized likelihood algorithm using the integral of the square of the first derivative in the penalty term, we have found by experience that the slightly more complicated second derivative approach gives less locally "rough" estimators. Namely we consider the problem

$$(4.1) \max L(f); f \in H_0^2(a, b), f \geq 0 \text{ and } \int_a^b f(x) dx = 1$$

where

$$(4.2) L(f) = \prod_{i=1}^n f(x_i) \exp[-\alpha \int_a^b |f''(x)|^2 dx], (\alpha > 0).$$

The details of the algorithm to approximate the solution to this problem are omitted, since they are very similar to the argument in Section 2. The operation count is still $O(n)$ per iteration.

In Figure 1, we demonstrate the solution to (4.1) using a random sample of size 20 from the standard normal distribution with mean 0 and variance 1. In Figure 2, we show the

estimator based on a sample of size 100. In Figures 3 and 4 we show the D.M.P.L.E. estimators for the 50-50 mixture of two normal distributions, both having variance 1 and with means at -1.5 and +1.5 for samples of size 25 and 100 respectively.

One comforting feature of the maximum penalized likelihood procedure is the relatively robust quality of the estimator in that changes of the optimal α with N and from distribution to distribution tend not to be traumatic, and that a rough and ready guess for α (e.g., 10) is frequently satisfactory. In Figures 5 and 6 we show an estimate for the Gaussian mixture mentioned above for a sample size of 300 and α values of 10 and .1.

If the density to be estimated is denoted by $f(\cdot)$ and the D.M.P.L.E. is denoted by $\hat{f}(\cdot)$, then we consider as one measure of estimate quality the average integrated mean square error

$$(4.3) \text{ IMSE} = \int (\hat{f}(x) - f(x))^2 f(x) dx .$$

I.M.S.E.'s for various α and N are given in Table 1 for the standard normal, the 50-50 normal mixture mentioned above, the t distribution with 5 degrees of freedom and the F distribution with (10,10) degrees of freedom.

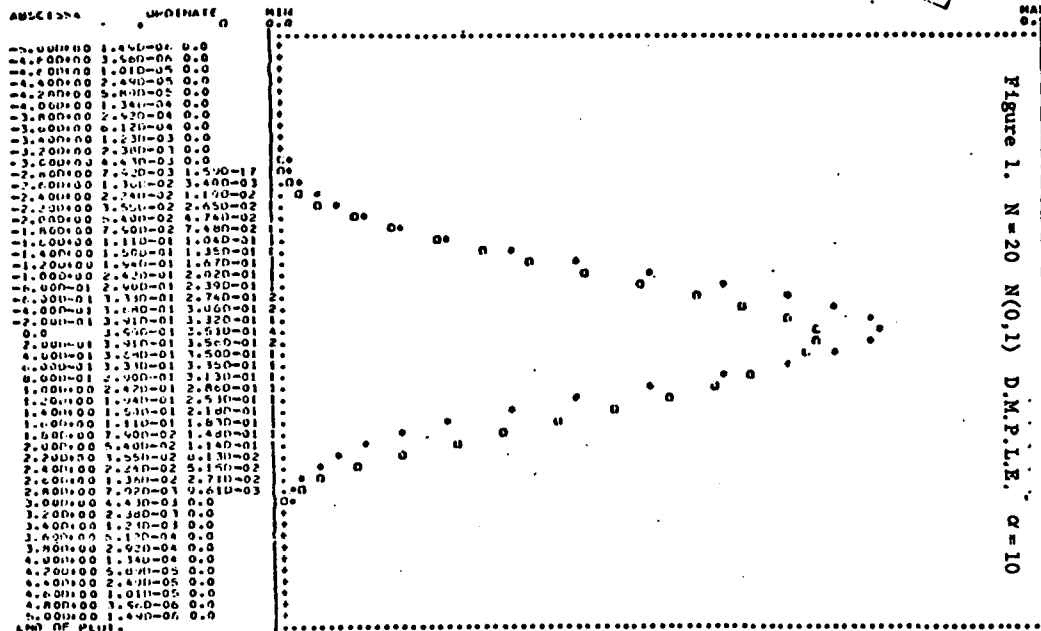
REFERENCES

1. Good, I.J. and Gaskins, R.A. (1971). "Nonparametric roughness penalties for probability densities," *Biometrika* 58, pp. 255-277.
2. de Montricher, G., Tapia, R.A., and Thompson, J.R. (1975). "Nonparametric maximum likelihood estimation of probability densities by penalty function methods," *Annals of Statistics* 3, pp. 1329-1348.
3. Tapia, R.A. (1974). "A stable approach to Newton's method for general mathematical programming problems in R^n ," *Journal of Optimization Theory and Applications* 14, pp. 453-476.
4. Tapia, R.A. (1976), "Diagonalized multiplier methods and quasi-Newton methods for constrained optimization," to appear in *Journal of Optimization Theory and Applications*.

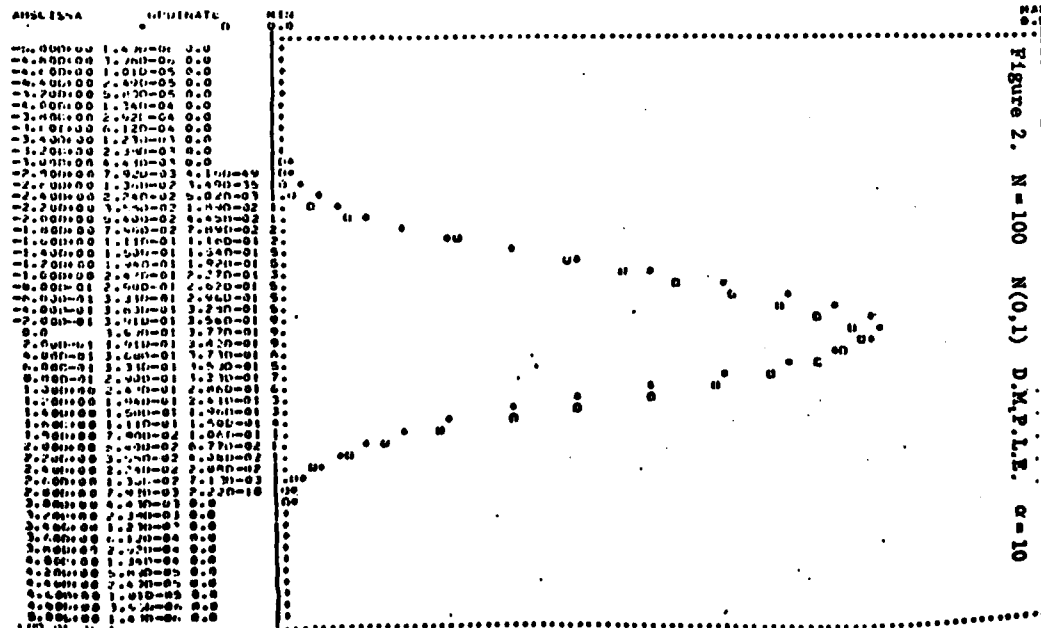
TABLE 1
Average I.M.S.E. of the D.M.P.L.E. for α Perturbed by a Factor of Two. Divide α by 10 for the $F_{10,10}$ Samples.

Sample	I.M.S.E. for		
	$\alpha = 5$	$\alpha = 10$	$\alpha = 20$
$N(0,1) N = 25$.00242	.00267	.00427
$N(0,1) N = 100$.00093	.00079	.00089
$N(0,1) N = 400$.00037	.00033	.00035
$N(0,1) N = 800$.00028	.00022	.00019
Bimodal $N = 25$.00197	.00159	.00152
Bimodal $N = 100$.00070	.00054	.00171
Bimodal $N = 400$.00030	.00024	.00022
$t_5 N = 25$.00297	.00282	.00350
$t_5 N = 100$.00092	.00084	.00101
$t_5 N = 400$.00039	.00032	.00030
$F_{10,10} N = 25$.03208	.03865	.05519
$F_{10,10} N = 100$.00996	.01390	.02411
$F_{10,10} N = 400$.00292	.00450	.00740

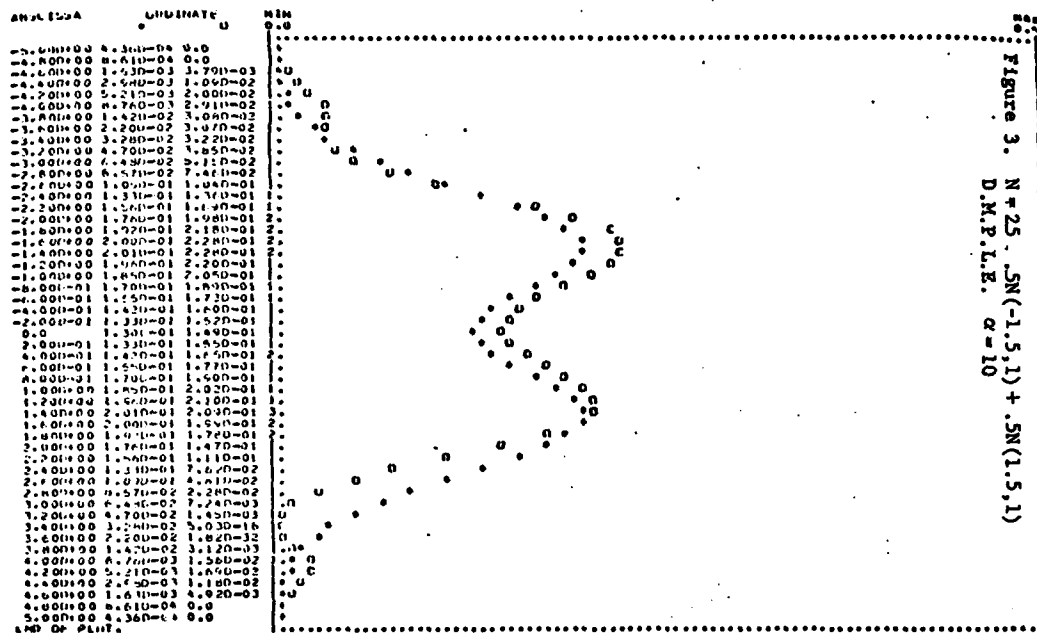
NORMAL 10.13 SAMPLE OF SIZE 70 VS.
 DISCRETIZED MAXIMUM LIKELIHOOD PENALIZED ESTIMATE
 WITH WEIGHTING PARAMETER ALPHA 0.1000000000000000
 27 MESH POINTS FROM -3.25000 TO 3.25000
 DISCRETIZED MESH INTERVAL 0.25000
 INTEGRATED MEAN SQUARE ERROR 0.28104018000000
 INTEGRATED SQUARE ERROR 0.45472018100000
 MAXIMUM ABSOLUTE DIFFERENCE 0.77015451300000
 LOG LIKELIHOOD TERM -0.27241002000000
 LOG PENALTY TERM -0.12554057070000



NORMAL 10.13 SAMPLE OF SIZE 100 VS.
 DISCRETIZED MAXIMUM LIKELIHOOD PENALIZED ESTIMATE
 WITH WEIGHTING PARAMETER ALPHA 0.1000000000000000
 27 MESH POINTS FROM -3.25000 TO 3.25000
 DISCRETIZED MESH INTERVAL 0.25000
 INTEGRATED MEAN SQUARE ERROR 0.79900000000000
 INTEGRATED SQUARE ERROR 0.34474554900000
 MAXIMUM ABSOLUTE DIFFERENCE 0.47101379100000
 LOG LIKELIHOOD TERM -1.11010425700000
 LOG PENALTY TERM -0.11013375400000



MINIMUM MINIMAL WITH MEANS = 1.50 VARIANCE OF LEFT = 1 WITH WEIGHT AND VARIANCE OF RIGHT = 0.5000 1.0000
 SAMPLE SIZE = 25 50 SAMPLES IN RIGHT VS.
 DISCRETE MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 WITH WEIGHTING PARAMETER ALPHA 0.10000000000000002
 41 MESH POINTS FROM -5.00000 TO 5.00000
 DISCRETE MESH INTERVAL 0.25000
 INTEGRATED MEAN SQUARE ERROR 0.70110744300-03
 INTEGRATED SQUARE ERROR 0.62201220670-02
 MAXIMUM ABSOLUTE DIFFERENCE 0.7347620520-01
 LOG LINE INHOD TERM -0.45213763870-02
 LOG PENALTY TERM -0.23516104510-01



MINIMUM MINIMAL WITH MEANS = 1.50 VARIANCE OF LEFT = 1 WITH WEIGHT AND VARIANCE OF RIGHT = 0.5000 1.0000
 SAMPLE SIZE = 100 50 SAMPLES IN RIGHT VS.
 DISCRETE MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 WITH WEIGHTING PARAMETER ALPHA 0.10000000000000002
 41 MESH POINTS FROM -5.00000 TO 5.00000
 DISCRETE MESH INTERVAL 0.25000
 INTEGRATED MEAN SQUARE ERROR 0.1124495640-01
 INTEGRATED SQUARE ERROR 0.5764400620-01
 MAXIMUM ABSOLUTE DIFFERENCE 0.23239250840-01
 LOG LINE INHOD TERM -0.18670774940-03
 LOG PENALTY TERM -0.17131462430-01

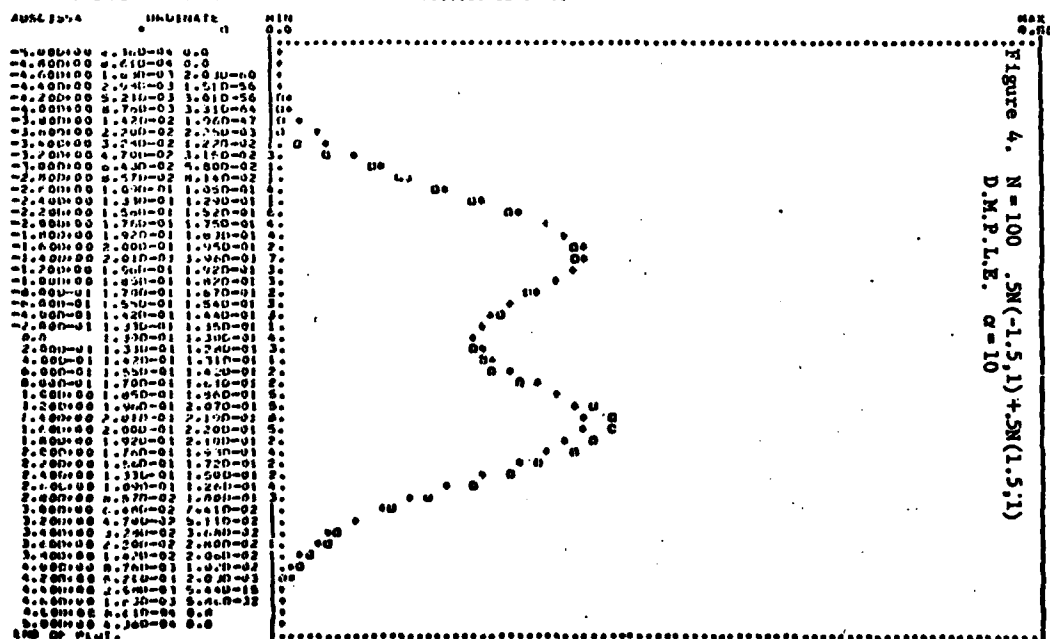


Figure 5. $N = 300$, $SN(-1.5, 1) + SN(1.5, 1)$
D.M.P.L.E., $\alpha = 10$

Figure 5. $N = 300$.5N(-1.5,1) + .5N(1.5,1)
D.M.P.L.E. $\alpha = 10$

Figure 6. $N=300$, $SM(-1.5,1) + .5M(1.5,1)$
D.M.P.L.E. $\alpha=.1$

Figure 6. $N = 300$. $.5N(-1.5, 1) + .5N(1.5, 1)$
D.M.P.L.E. $\alpha = .1$

END

DATE
FILMED

10-81

DTIC